DOCUMENT RESUME

ED 441 031 TM 030 830

AUTHOR Fahoome, Gail; Sawilowsky, Shlomo S.

TITLE Review of Twenty Nonparametric Statistics and Their Large

Sample Approximations.

PUB DATE 2000-04-00

NOTE 42p.; Paper presented at the Annual Meeting of the American

Educational Research Association (New Orleans, LA, April

24-28, 2000).

PUB TYPE Information Analyses (070) -- Numerical/Quantitative Data

(110) -- Speeches/Meeting Papers (150)

EDRS PRICE MF01/PC02 Plus Postage.

DESCRIPTORS Monte Carlo Methods; *Nonparametric Statistics; *Sample

Size; *Statistical Distributions

ABSTRACT

Nonparametric procedures are often more powerful than classical tests for real world data, which are rarely normally distributed. However, there are difficulties in using these tests. Computational formulas are scattered throughout the literature, and there is a lack of availability of tables of critical values. This paper brings together the computational formulas for 20 commonly used nonparametric tests that have large-sample approximations for the critical value. Because there is no generally agreed upon lower limit for the sample size, Monte Carlo methods have been used to determine the smallest sample size that can be used with the large-sample approximations. The statistics reviewed include single-population tests, comparisons of two populations, comparisons of several populations, and tests of association. (Contains 4 tables and 59 references.) (Author/SLD)



Review of Twenty Nonparametric Statistics and Their Large Sample Approximations

Gail Fahoome

Shlomo S. Sawilowsky

U.S. DEPARTMENT OF EDUCATION Office of Educational Research and Improvement EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)

This document has been reproduced as received from the person or organization originating it.

- ☐ Minor changes have been made to improve reproduction quality.
- Points of view or opinions stated in this document do not necessarily represent official OERI position or policy.

PERMISSION TO REPRODUCE AND DISSEMINATE THIS MATERIAL HAS BEEN GRANTED BY

0_, 11,00

TO THE EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)

AERA

BEST COPY AVAILABLE



Gail Fahoome

Mathematics Department

Detroit College of Business

Shlomo S. Sawilowsky

Educational Evaluation & Research

College of Education

Wayne State University

running head: Twenty Nonparametric Statistics

Gail Fahoome is adjunct instructor of mathematics at the Detroit College of Business, and adjunct instructor of Educational Evaluation and Research, College of Education, Wayne State University. Contact her at 18968 Hamburg, Detroit, MI 48205 for all communications regarding this paper. E-mail at gfahoome@aol.com. Her areas of expertise are mathematics education, Monte Carlo methods, nonparametric statistics, and program evaluation.

Shlomo Sawilowsky is professor and chair of Educational Evaluation and Research, 351 College of Education, Wayne State University, Detroit, MI 48202. E-mail at shlomo@edstat.coe.wayne.edu. His areas of expertise are Monte Carlo methods, rank statistics, research and experimental design, and classical measurement.

The following is an historical note regarding some of the authors cited in this paper, with the date of doctoral dissertation in parenthesis. R. Clifford Blair was the major professor to Shlomo S. Sawilowsky (1985; James J. Higgins was a member of the dissertation committee) and Theodore Micceri (1986) at the University of South Florida. Shlomo S. Sawilowsky was the major professor to D. Lynn Kelley (1994), Patrick D. Bridge (1996), Margaret P. Posch (1996), Todd C. Headrick (1997), Michael J. Nanna (1997), Joseph Musial, III (1999), and Gail F. Fahoome (1999) at Wayne State University.



Nonparametric procedures are often more powerful than classical tests for real world data which are rarely normally distributed. However, there are difficulties in using these tests. Computational formulas are scattered throughout the literature, and there is a lack of availability of tables of critical values. We bring together the computational formulas for twenty commonly employed nonparametric tests that have large-sample approximations for the critical value. Because there is no generally agreed upon lower limit for the sample size, we use Monte Carlo methods to determine the smallest sample size that can be used with the large-sample approximations. The statistics reviewed include single-population tests, comparisons of two populations, comparisons of several populations, and tests of association.



Classical parametric tests, such as the F and t, were developed in the early part of the twentieth century. These statistics require the assumption of population normality. Bradley (1968) wrote, "To the layman unable to follow the derivation but ambitious enough to read the words, it sounded as if the mathematician had esoteric mathematical reasons for believing in at least quasi-universal quasi-normality" (p. 8). "Indeed, in some quarters the normal distribution seems to have been regarded as embodying metaphysical and awe-inspiring properties suggestive of Divine Intervention" (p. 5).

However, when Micceri (1989) investigated 440 large-sample education and psychology data sets, he concluded "No distributions among those investigated passed all tests of normality, and very few seem to be even reasonably close approximations to the Gaussian" (p. 161). This is of great practical importance because even though the well known Student's *t* test is preferable to nonparametric competitors when the normality assumption has been met, Blair and Higgins (1980) noted:

Generally unrecognized, or at least not made apparent to the reader, is the fact that the t test's claim to power superiority rests on certain optimal power properties that are obtained under normal theory. Thus, when the shape of the sampled population(s) is unspecified, there are no mathematical or statistical imperatives to ensure the power superiority of this statistic. (p. 311)

Blair and Higgins (1980) demonstrated the power superiority of the nonparametric Wilcoxon Rank Sum test over the *t* test for a variety of nonnormal theoretical distributions. In a Monte Carlo study of Micceri's real world data sets, Sawilowsky and Blair (1992) concluded that although the *t* test is generally robust with respect to Type I errors under conditions of equal sample size, fairly large samples, and two-tailed tests, it is not powerful for skewed distributions. Under these conditions, the Wilcoxon Rank Sum test is three to four times more powerful. See also Bridge and Sawilowsky (1999) and Nanna and Sawilowsky (1998).

It is appropriate to consider further this class of statistics because of the power advantages of nonparametric tests with real world data. The terms 'nonparametric' and 'distribution-free' are often used interchangeably to describe tests that make few, if any, assumptions about the distribution of the population. There is, however, a distinction between them. Bradley (1968) explained that "a



nonparametric test is one which makes no hypothesis about the value of a parameter in a statistical density function, whereas a distribution-free test is one which makes no assumptions about the precise form of the sampled population" (p. 15). In this paper we are concerned with nonparametric procedures.

A difficulty in using nonparametric tests is the availability of computational formulas and tables of critical values. For example, Siegel and Castellan (1988) noted, "Valuable as these sources are, they have typically either been highly selective in the techniques presented or have not included the tables of significance" (p. xvi). This continues to be a problem as evidenced by our survey of 20 in-print generic college statistics textbooks, including seven general textbooks, eight for the social and behavioral sciences, four for business, and one for engineering. Formulas were given for only eight nonparametric statistics, and tables of critical values were given for only the following six: (a) Kolmogorov-Smirnov test, (b) Sign test, (c) Wilcoxon Signed Rank test, (d) Wilcoxon (Mann-Whitney) test, (e) Spearman's rank correlation coefficient, and (f) Kendall's rank correlation coefficient.

This situation is somewhat improved for nonparametric statistics textbooks. Eighteen nonparametric textbooks published since 1956 were also reviewed. The most comprehensive texts in terms of coverage were Neave and Worthington (1988) which is out of print and Deshpande Gore, and Shanubhogue (1995). Table 1 contains the statistical content of the eighteen textbooks. The comment by Laubscher, Steffens, and De Lange (1968) on the Mood test summarized the findings: "As far as we know the main drawback in using this test statistic, developed more than 14 years ago, lies in the fact that its distribution has never been tabulated except for a few isolated cases" (p. 497).

Table 1. Results of Survey of 18 Nonparametric Books.

Statistic	Number of Books That Included Tables of Critical Values		
Single Population Tests			
Kolgomorov-Smirnov Goodness-of-Fit Test	11		
Sign Test	4		
Wilcoxon's Signed Rank Test	14		
Comparison of Two Populations			
Kolmogorov-Smirnov Two Sample Test	11		
Rosenbaum's Test	1		



Wilcoxon (Mann-Whitney) Test	14
Mood Test	1
Savage Test	1
Ansari-Bradley Test	1
Comparison of Several Populations	
Kruskal-Wallis Test	10
Friedman's Test	9
Terpstra-Jonckheere Test	5
Page's Test	4
Match Test for Ordered Alternatives	1
Tests of Association	
Spearman's Rank Correlation Coefficient	12
Kendall's Rank Correlation Coefficient	10

Many nonparametric tests have large sample approximations that can be used as an alternative to tabulated critical values. These approximations are useful substitutes if the sample size is sufficiently large, and hence, obviate the need for locating tables of critical values. However, there is no generally agreed upon definition of what constitutes a *large* sample size. Consider the Sign test and the Wilcoxon tests as examples.

Regarding the Sign test, Hájek (1969) wrote, "The normal approximation is good for $N \ge 12$ " (p. 108). Gibbons (1971) agreed, "Therefore, for moderate and large values of N (say at least 12) it is satisfactory to use the normal approximation to the binomial to determine the rejection region" (p. 102). Both Sprent (1989) and Deshpande, Gore, and Shanubhogue (1995), however, recommended n greater than 20. Siegel and Castellan (1988) suggested $n \ge 35$, but Neave and Worthington (1988) proposed that n > 50.

The literature regarding the Wilcoxon Rank Sum test is similarly disparate. Deshpande, Gore, and Shanubhogue (1995) stated that the combined sample size should be at least 20 to use a large sample approximation of the critical value. Conover (1971) and Sprent (1989) recommended that one or both samples must exceed 20. Gibbons (1971) placed the lower limit at twelve per sample. For the Wilcoxon Signed Rank test, Deshpande, Gore, and Shanubhogue (1995) said that the approximation can be used when n is greater than 10. Gibbons (1971) recommended it when n is greater than 12, and Sprent (1989)



required n to be greater than 20. The general lack of agreement may indicate that these recommendations are based on personal experience, the sample sizes in available tables, the author's definition of "acceptable" or "large", or some other criterion.

There are two alternatives to tables and approximations. The first is to use exact permutation methods. There is software available that will generate exact p-values for *small* data sets and Monte Carlo estimates for *larger* problems. See Ludbrook and Dudley (1998) for a brief review of the capabilities of currently available software packages for permutation tests. However, these software solutions are expensive, have different limitations in coverage of procedures, and may require considerable computing time even with fast personal computers (see, e.g., Musial, 1999; Posch & Sawilowsky, 1997). In any case, a desirable feature of nonparametric statistics is that they are easy to compute without statistical software and computers, which makes their use in the classroom or work in the field attractive.

A second alternative is the use of the rank transformation (RT) procedure developed by Conover and Iman (1981). They proposed the use of this procedure as a bridge between parametric and nonparametric techniques. The RT is carried out as follows: rank the original scores, perform the classical test on the ranks, and refer to the standard table of critical values. In some cases, this procedure results in a well-known test. For example, conducting the *t* test on the ranks of original scores in a two independent samples layout is equivalent to the Wilcoxon Rank Sum test. (However, see the caution noted by Sawilowsky & Brown, 1991). In other cases, such as factorial analysis of variance (ANOVA) layouts, a new statistic emerges.

The early exuberance with this procedure was related to its simplicity and promise of increased statistical power when data sets displayed nonnormality. Iman and Conover noted the success of the RT in the two independent samples case and the one-way ANOVA layout. Nanna (1997) showed that the RT is robust and powerful as an alternative to the independent samples multivariate Hotelling's T².

However, Blair and Higgins (1985) demonstrated that the RT suffers power losses in the dependent samples t test layout as the correlation between the pretest and posttest increases. Bradstreet (1997) found the RT to perform poorly for the two sample Behrens-Fisher problem. Sawilowsky (1985), Sawilowsky, Blair, and Higgins (1989), Blair, Sawilowsky, and Higgins (1987), and Kelley and Sawilowsky (1997) showed the RT has severely inflated Type I errors and a lack of power in testing interactions in factorial ANOVA layouts. Harwell and Serlin (1997) found the RT to have inflated Type I errors in the test of $\beta = 0$ in linear regression. In the context of analysis of covariance, Headrick and

. 1



Sawilowsky (1999, 2000) found the RT's Type I error rate inflates quicker than the general ANOVA case, and it demonstrated more severely depressed power properties. Recent results by Headrick (personal communications) shows the RT to have poor control of Type I errors in the ordinary least squares multiple regression layout. Sawilowsky (1989) stated that the RT as a bridge has fallen down, and cannot be used to unify parametric and nonparametric methodology or as a method to avoid finding formulas and critical values for nonparametric tests.

The Current Study

As noted above, the computational formulas for many nonparametric tests are scattered throughout the literature, and tables of critical values are scarcer. Large sample approximation formulas are also scattered and appear in different forms. Most important, the advice on how "large" a sample must be to use the approximations is conflicting. The purpose of this study is to ameliorate all five of these problems.

Ascertaining the smallest sample size that can be used with a large sample approximation for the various statistics would enable researchers who do not have access to the necessary tables of critical values or statistical software to employ these tests. The first portion of this paper uses Monte Carlo methods to determine the smallest sample size that can be used with the large sample approximation while still preserving nominal alpha. The second portion of this paper provides a comprehensive review of computational formulas with worked examples for twenty nonparametric statistics. They were chosen because they are commonly employed and because large sample approximation formulas have been developed for them.

Methodology

Each of the twenty statistics was tested with normal data and Micceri's (1989; see also Sawilowsky, Blair, & Micceri, 1990) real world data sets. The real data sets represent smooth symmetric, extreme asymmetric, and multi-modal lumpy distributions. Monte Carlo methods were used in order to determine the smallest samples that can be used with large-sample approximations.

A program was written in Fortran 90 (Lahey, 1998) for each statistic. The program sampled with replacement from each of the four data sets for n = 1, 2, ..., N; $n_1 = n_2 = (2, 2), (3,3), ..., (N_1,N_2)$, and so forth as the number of groups increased. The statistic was calculated and evaluated using the tabled values when available and the approximation of the critical value. The number of rejections was counted



and the Type I error rate was computed. Nominal α was set at .05 and .01. Bradley's (1978) conservative estimates of .045 < Type I error rate < .055 and .009 < Type I error rate < .011 were used, respectively, as measures of robustness. The sample sizes were increased until the Type I error rates converged within these acceptable regions.

Assumptions and Limitations

In many cases there are different formulas for the large sample approximation of a statistic. Two criteria were used in choosing which formula to include: (a) consensus of authors, and (b) ease of use in computing and programming. Some of the statistics have different large sample approximations based on the presence of ties among the data. The formulas not based on ties were used because we corrected for ties using average ranks.

Data Sets For Worked Examples In This Article

The worked examples in this study used five data sets that may be found in Table 3 (Appendix). Some statistics converged at relatively large sample sizes. In choosing the sample size for the worked example, we compromised between the amount of computation required for large samples and an unrepresentatively small but convenient sample size. Therefore, we selected a sample size of n = 15, recognizing that some statistics' large sample approximations do not converge within Bradley's (1968) limits for this small sample size. The data sets were randomly selected from Micceri's (1989) multimodal lumpy data set, Table 4 (Appendix). Because the samples came from the same population, the worked examples all conclude that the null hypothesis cannot be rejected.

Statistics Examined

The twenty statistics included in this article represent four layouts: (1) single population tests, (2) comparison of two populations, (3) comparison of several populations, and (4) tests of association. Single-populations tests included: (a) a goodness-of-fit test, (b) tests for location, and (c) an estimator of the median. Comparisons of two populations included: (a) tests for general differences, (b) two-sample location problems, and (c) two-sample scale problems. Comparisons of several populations included: (a) ordered alternative hypotheses, and (b) tests of homogeneity against omnibus alternatives. Tests of association focused on rank correlation coefficients.



BEST COPY AVAILABLE

STATE OF THE STATE

Results

Table 2 shows the minimum sample sizes for the tests studied. These recommendations are based on results that converged when underlying assumptions are reasonably met. The minimum sample-sizes are conservative, representing the largest minimum for each test. If the test had three or more samples, the largest group minimum was chosen. Consequently the large-sample approximations will work in some instances for smaller sample sizes. Where the test involves more than one sample, the smallest sample size refers to the smallest sample size for each equal sample.

Table 2. Minimum Sample Size for Large-Sample Approximations.

Test	α = .05	α =.01
Single Population Tests	_	
Kolmogorov-Smirnov Goodness-of-Fit Test	$25 \le n \le 40$	$28 \le n \le 50$
Sign Test	n > 150	n > 150
Wilcoxon Signed Rank Test	10	22
Estimator of Median for a Continuous Distribution	n > 150	n > 150
Comparison of Two Populations		
Kolmogorov-Smirnov Two Sample Test	n > 150	n > 150
Rosenbaum's Test	16	20
Tukey's Test	$10 \le n \le 18$	21
Wilcoxon (Mann-Whitney) Test	15	29
Hodges-Lehmann Estimator	15	20
Siegel-Tukey Test	25	38
Mood Test	5	23
Savage Test	11	31
Ansari-Bradley Test	16	29
Comparison of Several Populations		
Kruskal-Wallis Test	11	22
Friedman's Test	13	23
Terpstra-Jonckheere Test	4 .	8
Page's Test $(k > 4)$	11	18
The Match Test for Ordered Alternatives $(k > 3)$	86	27



<u>Tests of Association</u>
Spearman's Rank Correlation Coefficient
Kendall's Rank Correlation Coefficient

12 40 $14 \le n \le 24$ $15 \le n \le 35$

Some notes and cautionary statements are in order with regard to the entries in Table 3. The Monte Carlo methods were completed for n = 1, 2, ... 150. The Kolmogorov-Smirnov goodness-of-fit test was conservative for values below the minimum value stated and liberal for values above the maximum value. Results for the Sign test indicate convergence for some distributions may occur close to n = 150. The results for the confidence interval for the Estimator of the Median suggest convergence may occur close to n = 150 only for normally distributed data. However, for the nonnormal data sets the Type I error rates were quite conservative (e.g., for $\alpha = .05$ the Type I error rate was only 0.01146 and for $\alpha = .01$ it was only 0.00291 for n = 150 and the extreme asymmetric data set).

The Kolmogorov-Smirnov test was erratic, with no indication convergence would be close to 150. Results for Tukey's Test were conservative for $\alpha=.05$ when the cutoff for the p-value was .05, and fell within acceptable limits for some sample sizes when .055 was used as a cutoff. The Hodges-Lehmann Estimator only converged for normal data. For nonnormal data the large sample approximation was extremely conservative with n=10 (e.g., for the extreme asymmetric data set the Type I error rate was only 0.0211 and 0.0028 for the .05 and .01 alpha levels, respectively) and increased in conservativeness (i.e., the Type I error rate converged to 0.0) as n increased. The Match test only converged for normally distributed data, and it was the only test where the sample size required for $\alpha=.01$ was smaller than for $\alpha=.05$.

Statistics, Worked Examples, Large Scale Approximations <u>Single Population Tests</u>

Goodness-of-fit statistics are single-population tests of how well observed data fit expected probabilities or a theoretical probability density function. They are often used as a preliminary test of the distribution assumption of parametric tests. The Kolmogorov-Smirnov goodness-of-fit test was studied.

Tests for location are used to make inferences about the location of a population. The measure of location is usually the median. If the median is not known but there is reason to believe that its value is M_0 , then the null hypothesis is $H_0: M = M_0$. The tests for location studied were the Sign test, Wilcoxon's Rank Sum test, and the Estimator of the Median for a Continuous Distribution.



Kolmogorov-Smirnov Goodness-of-Fit Test

The Kolmogorov-Smirnov (K-S) goodness-of-fit statistic was devised by Kolmogorov in 1933 and Smirnov in 1939. It is a test of goodness-of-fit for continuous data, based on the maximum vertical deviation between the empirical distribution function, $F_n(x)$, and the hypothesized cumulative distribution function, $F_0(x)$. Small differences support the null hypothesis while large differences are evidence against the null hypothesis.

The null hypothesis is $H_0: F_n(x) = F_0(x)$ for all x and the alternative hypothesis is $H_1: F_n(x) \neq F_0(x)$ for at least some x where $F_0(x)$ is a completely specified continuous distribution. The empirical distribution function, $F_n(x)$, is a step function, defined as:

$$F_n(x) = \frac{\text{number of sample values} \le x}{n} \tag{1}$$

where n = sample size.

Test statistic.

The test statistic, D_n , is the maximum vertical distance between the empirical distribution function and the cumulative distribution function.

$$D_n = \max \left[\max \left| F_n(x_i) - F_0(x_i) \right|, \max \left| F_n(x_{i-1}) - F_0(x_i) \right| \right]$$
 (2)

Both vertical distances $F_n(x_i) - F_0(x_i)$ and $F_n(x_{i-1}) - F_0(x_i)$ have to be calculated in order to find the maximum deviation. The overall maximum of the two calculated deviations is defined as D_n .

For a one-tailed test against the alternatives $H_1: F_n(x) > F_0(x)$ or $H_1: F_n(x) < F_0(x)$ for at least some values of x, the test statistics are respectively:

$$D_n^* = \max[F_n(x) - F_0(x)]$$
 (3)

or

$$D_n^- = \max[F_0(x) - F_n(x)]$$
 (4)

The rejection rule is to reject H_0 when $D_n \ge D_{n,\alpha}$ where $D_{n,\alpha}$ is the critical value for a given n and α level of significance.

Large sample sizes.

The null distribution of $4nD_n^{+2}$ (or $4nD_n^{-2}$) is approximately χ^2 with 2 degrees of freedom. Thus, the large sample approximation is

BEST COPY AVAILABLE



$$D_n^+ \approx \frac{\sqrt{\chi_{\alpha,2}}}{2\sqrt{n}} \approx \frac{1}{2} \sqrt{\frac{\chi_{\alpha,2}^2}{n}} \tag{5}$$

where $\chi_{\alpha,2}^2$ is the value for chi-square with 2 degrees of freedom for the appropriate alpha level and n is the sample size.

Example.

The K-S goodness-of-fit statistic was calculated for Sample 1 in Table 3 (Appendix), n = 15, against the cumulative frequency distribution of the multimodal lumpy data set. The maximum difference at step was 0.07463 and the maximum difference before step was 0.142610. Thus the value of D_n is 0.142610. For a two-tail test with $\alpha = .05$, the large sample approximation is $1.3581/\sqrt{n} = 1.3581/\sqrt{15} = 0.35066$. Because 0.142610 < 0.35066, the null hypothesis cannot be rejected.

The Sign Test

The Sign test is credited to Fisher as early as 1925. One of the first papers on the theory and application of the sign test is attributed to Dixon and Mood in 1946 (Hollander & Wolfe, 1973). According to Neave and Worthington (1988), the logic of the Sign test is "almost certainly the oldest of all formal statistical tests as there is published evidence of its use long ago by J. Arbuthnott (1710)!" (p. 65).

The Sign test is a test for a population median. It can also be used with matched data as a test for equality of medians. The test is based upon the number of values above or below the hypothesized median. Gibbons (1971) referred to the sign test as the nonparametric counterpart of the one-sample t test. The sign test tests the null hypothesis $H_0: M = M_0$ where M is the sample median and M_0 is the hypothesized population median against the alternative hypothesis $H_1: M \neq M_0$. One-tailed test alternative hypotheses are of the form $H_1: M < M_0$ and $H_1: M > M_0$.

Procedure.

Each x_i is compared with M_0 . If $x_i > M_0$ then a plus sign '+' is recorded. If $x_i < M_0$ then a minus sign '-' is recorded. In this way all data are reduced to '+' and '-' signs.

Test statistic.

The test statistic is the number of '+' signs or the number of '-' signs. If the expectation under the alternative hypothesis is that there will be a preponderance of '+' signs, the test statistic is the



BEST COPY AVAILABLE

number of '-' signs. Similarly, if the expectation is a preponderance of '-' signs, the test statistic is the number of '+' signs. If the test is two-tailed, use the smaller of the two. Thus,

$$S =$$
the number of '+' or '-' signs (depending upon the context) (6)

Large sample sizes.

The large sample approximation is given by

$$S^{\bullet} = \frac{S - \frac{n}{2}}{\sqrt{\frac{n}{4}}} \tag{7}$$

where S is the test statistic and n is the sample size. S^* is compared to the standard normal z scores for the appropriate α level.

Example.

The Sign test was calculated using Sample 1 in Table 3 (Appendix), n = 15. The population median is 18.0. The number of negative values is 7 and the number of positive values is 8. Therefore S = 7. The large sample approximation, S^* , using formula (7) is -.258199. Because -.258199 > -1.95996, the null hypothesis cannot be rejected.

Wilcoxon's Signed Rank Test

Wilcoxon's Signed Rank test was introduced by Wilcoxon in 1945. The statistic uses the ranks of the absolute differences between x_i and M_0 along with the sign of the difference. This uses the relative magnitudes of the data. This statistic can also be used to test for symmetry and to test for equality of location for paired replicates.

The null hypothesis is $H_0: M = M_0$ against the alternative $H_1: M \neq M_0$. The alternative may also be one-sided, $H_1: M > M_0$ or $H_1: M < M_0$.

Procedure.

Compute the differences, D_i , by the formula

$$D_i = x_i - M_0 \tag{8}$$

Rank the absolute value of the differences, in ascending order, keeping track of the individual signs.



Test statistic.

The test statistic is the sum of either the positive ranks or the negative ranks. If the alternative hypothesis suggests that the sum of the positive ranks should be large,

then

$$T$$
 = the sum of ranks of the negative differences (9)

If the alternative hypothesis suggests that the sum of the negative ranks should be large, then

$$T^{+}$$
 = the sum of ranks of the positive differences (10)

For a two-tailed test, T is the smaller of the two rank-sums. The total sum of the ranks is $\frac{n(n+1)}{2}$, which gives the following relationship:

$$T^{+} = \frac{n(n+1)}{2} - T^{-} \tag{11}$$

Large sample sizes.

The large sample approximation is

$$z = \frac{T - \frac{n(n+1)}{4}}{\sqrt{\frac{n(n+1)(2n+1)}{24}}}$$
(12)

where T is the test statistic and n is the sample size. The resulting z is compared to the standard normal z for the appropriate alpha level.

Example.

The Signed Rank test was computed using the data from Sample 1 in Table 3 (Appendix), n = 15. The median of the population is 18.0. Tied differences were assigned midranks. The sum of the negative ranks was 38.5 and the sum of the positive ranks was 81.5. Therefore the Signed Rank statistic is 38.5. The large sample approximation is $\frac{-21.5}{\sqrt{310}} = \frac{-21.5}{17.6068} = -1.22112$. Because -1.22112 > -1.95996, the null hypothesis is not rejected.

Estimator of the Median for a Continuous Distribution

The sample median is the point estimate of the population median. This procedure provides a $1-\alpha$ confidence interval for the population median. It was designed to be used with continuous data.



Procedure.

Let n be the size of the sample. Order the n observations in ascending order, $x_{(1)} \le x_{(2)} \le \ldots \le x_{(n)}$. Let $x_{(0)} = -\infty$ and $x_{(n+1)} = \infty$. These n+2 values form n+1 intervals $(x_{(0)}, x_{(1)}), (x_{(1)}, x_{(2)}), \ldots, (x_{(n-1)}, x_{(n)}), (x_{(n)}, x_{(n+1)})$. The ith interval is defined as $(x_{(i-1)}, x_{(i)})$ with $i=1, 2, \ldots, n, n+1$. The probability that the median is in any one interval is based on the binomial distribution. The confidence interval for the median given the confidence coefficient $1-\alpha$, requires that an r be found such that the sum of the probabilities of the intervals in both the lower and upper ends give the best conservative approximation of $\alpha/2$, according to the following:

$$\sum_{j=0}^{r} \binom{n}{j} \frac{1}{2^n} \approx \frac{\alpha}{2} \approx \sum_{j=n-r}^{n} \binom{n}{j} \frac{1}{2^n}$$
 (13)

Thus $(x_{(r)}, x_{(r+1)})$ is the last interval in the lower end making $x_{(r+1)}$ the lower limit of the confidence interval. By a similar process, $x_{(n-r)}$ is the upper limit of the confidence interval.

Large sample sizes.

According to Deshpande, Gore, and Shanubhogue (1995) "one may use the critical points of the standard normal distribution, to choose the value of r + 1 and n - r, in the following way": r + 1 is the integer closest to

$$\frac{n}{2} - z_{\alpha/2} \left(\frac{n}{4}\right)^{\frac{1}{2}} \tag{14}$$

where $z_{\alpha/2}$ is the upper $\alpha/2$ critical value of the standard normal distribution.

Example.

The data from Sample 1 in Table 3 (Appendix), n = 15, were used to compute the estimator of the median. The population median is 18.0. For the given n and $\alpha = .05$, the value of r is 3. The value of r + 1 is 4, and n - r is 12. The 4th value is 13 and the 12^{th} value is 33. Therefore the interval is (13, 33). The large sample approximation yields 7.5 - 1.95996(1.9365) = 7.5 - 3.70 = 3.80. The closest integer is r + 1 = 4, so r = 3 and n - r = 12, resulting in the same interval, (13, 33). The interval contains the population median, 18.



Two Sample Problems

The two-sample problem consists of two independent random samples drawn from two populations. This study examined two sample tests for general differences, two sample location problems, and two sample scale problems.

When differences between two samples are not expected to be predominantly differences in location or differences in scale, a test for general differences is appropriate. Generally differences in variability are related to differences in location. Two tests for differences were considered, the Kolmogorov-Smirnov test for general differences and Rosenbaum's test.

Two sample location problems involve tests for a difference in location between two samples when the populations are assumed to be similar in shape. The idea is that $f_1(x) = f_2(x+\theta)$ or $f_1(x) = f_2(x-\theta)$ where θ is the distance between the population medians. Tukey's quick test, the Wilcoxon (Mann-Whitney) statistic, and the Hodges-Lehmann estimator of the difference in location for two populations were considered.

In two sample scale problems, the population distributions are usually assumed to have the same location with different spreads. However, Neave and Worthington (1988) cautioned that tests for difference in scale could be severely impaired if there is a difference in location as well. The following nonparametric tests for scale were studied: the Siegel-Tukey test, the Mood test, the Savage test for positive random variables, and the Ansari-Bradley test.

Kolmogorov-Smirnov Test for General Differences

The Kolmogorov-Smirnov test compares the cumulative distribution frequencies of the two samples to test for general differences between the populations of the samples. The sample cdf "is an approximation of the true cdf of the corresponding population – though, admittedly, a rather crude one if the sample size is small" (Neave & Worthington, 1988, p. 149). This property was used in the goodness-of-fit test above. Large differences in the sample cdf's can indicate a difference in the population cdf's, which could be due to differences in location, spread, or more general differences in the distributions. The null hypothesis is $H_0: F_1(x) = F_2(x)$ for all x and the alternative hypothesis is $H_1: F_1(x) \neq F_2(x)$ for some x.



Procedure.

The combined observations are ordered from smallest to largest, keeping track of the sample membership. Above each score, write the cdf of sample 1, and below each score write the cdf of sample 2. Because the samples are of equal sizes, it is only necessary to use the numerator of the cdf. For example, the $cdf(x_i) = \frac{i}{n}$. Then write i above x_i for sample 1. Find the largest difference between the cdf for sample 1 and the cdf for sample 2.

Test statistic.

The test statistic is D^* . $D^* = mnD$, and $D^* = n^2D$ for equal sample size. The above procedure yields nD. Thus

$$D^* = n(nD) \tag{15}$$

The greatest difference found by the procedure is multiplied by the sample size.

Large sample sizes.

As sample size increases, the distribution is approximately chi-squared with 2 degrees of freedom, as it is for the goodness-of-fit test. The large sample approximation for D is

$$D = \frac{1}{2} \sqrt{\frac{\chi_{\alpha,2}^2(m+n)}{mn}} \tag{16}$$

where $\chi_{a,2}^2$ is the value for chi-square with 2 degrees of freedom for the appropriate alpha level and n, m are the two sample sizes. The resulting D is used in formula (15).

Example.

This example used the data from Sample 1 and Sample 5 in Table 3 (Appendix), n = m = 15. The greatest difference (nD) between the cdf's of the two samples is nD = 3. Therefore $D^{\bullet} = 15(3) = 45$. The large sample approximation is $15^{2}(1.3581)\sqrt{\frac{30}{225}} = 225(1.3581)(.365148) = 111.579301$. Because 45 < 111.579301, the null hypothesis cannot be rejected.

Rosenbaum's Test

Rosenbaum's test, which was developed in 1965, is useful in situations where an increase in the measure of location implies an increase in variation. It is a quick and easy test based on the number of observations in one sample greater than the largest observation in the other sample.

The null hypothesis is that both populations have the same location and spread against the



alternative, that both populations differ in location and spread.

Procedure.

The largest observation in each sample is identified. If the largest overall observation is from sample 1, then the number of observations from sample 1 which are greater than the largest observation from sample 2 are counted. If the largest overall observation is from sample 2, then the number of observations from sample 2 which are greater than the largest observation from sample 1 are counted.

Test statistic.

The test statistic is the count of the extreme observations. R is the number of observations from sample 1 greater than the largest observation in sample 2 or the number of observations from sample 2 greater than the largest observation in sample 1.

Large sample sizes.

As sample sizes increase, $\frac{n_1}{N} \to p$ and the probability that the number of extreme values equals h approaches p^h .

Example.

Rosenbaum's statistic was calculated using Samples 1 and 5 in Table 3 (Appendix), $n_1 = n_2 = 15$. The maximum value from Sample 1 is 39, and from Sample 2, 33. There are three values from Sample 1 greater than 33, namely 34, 36, and 39. Hence R = 3. The large sample approximation is $(.5)^3 = 0.125$. Because 0.125 > .05, the null hypothesis cannot be rejected.

Tukey's Quick Test

Tukey published a quick and easy test for the two sample location problem in 1959. It is easy to calculate and in most cases does not require the use of tables. The most common one-tailed critical values are 6 ($\alpha = .05$) and 9 ($\alpha = .01$) for most sample sizes. The statistic is based on the sum of the extreme runs. If there is a difference in location between samples X and Y, one would expect more X's at one end and Y's at the other end when the combined samples are ordered.

Procedure.

The combined samples can be ordered, but it is only necessary to order the largest and smallest elements. If both the maximum and minimum value come from the same sample the test is finished, the value of $T_y = 0$, and the null hypothesis is not rejected.

For the one-tailed test, the lower end run should come from the sample expected to have the



lower median and the upper run from the sample expected to have the larger median. For a two-tailed test, it is possible to proceed with the test as long as the maximum and minimum come from different samples.

Test statistic.

 T_y is defined as follows for $H_1 = M_y > M_x$. T_y is the number of X's less than the smallest value of Y plus the number of Y's greater than the largest value of X. If $H_1 = M_x > M_y$ then the samples are reversed. For the two-tailed hypothesis both possibilities are considered.

Critical values.

As stated above, generally, the critical value for $\alpha = .05$ is 6, and is 9 for $\alpha = .01$. There are tables available. As long as the ratio of n_x to n_y is within 1 to 1.5, these critical values work well. There are corrections available when the ratio exceeds 1.5. For a two-tailed test the critical values are 7 ($\alpha = .05$) and 10 ($\alpha = .01$).

Large sample sizes.

The null distribution is based on the order of the elements of both samples at the extreme ends. It does not depend upon the order of the elements in the middle. The formula for the probability that $T_y \ge h$ is the sum of a finite geometric series,

$$\operatorname{Prob}(T_{y} \ge h) = \frac{pq(q^{h} - p^{h})}{q - p} \tag{17}$$

When the sample sizes are equal, p = q = .5. Then the probability of $T_y \ge h$ is $h \cdot 2^{-(h+1)}$. For a two-tailed test the probability is doubled.

Example.

The Tukey test was calculated using the data in Sample 1 and Sample 5 in Table 3 (Appendix), n = m = 15. The maximum value, 39, is from Sample 1 and the minimum, 2, is from Sample 5 so the test may proceed. The value of $T_y = 1 + 3 = 4$. For a two-tailed test with $\alpha = .05$, the large sample approximation is $2(4)(2^{-5}) = 0.25$. Because 0.25 > .05, the null hypothesis cannot be rejected.

Wilcoxon (Mann-Whitney) Statistic

In 1945, Wilcoxon introduced the Rank Sum test at the same time as the Signed Rank test. Mann and Whitney introduced a different version of the test in 1947. The Wilcoxon statistic is easily converted to the Mann-Whitney U statistic. The hypotheses of the test are $H_0: F_1(x) = F_2(x)$ for all x against the



two-tailed alternative, $H_0: F_1(x) \neq F_2(x)$. The one-tailed alternative is $H_1: F_1(x) = F_2(x+\theta)$.

Procedure.

For the Wilcoxon test, the combined samples are ordered, keeping track of sample membership. The ranks of the sample that is expected, under the alternative hypothesis, to have the smallest sum, are added. The Mann-Whitney test is as follows. Put all the observations in order, noting sample membership. Count how many of the observations of one sample exceed each observation in the first sample. The sum of these counts is the test statistic, U.

Test statistic.

For the Wilcoxon test,

$$S_n = \sum_{j=1}^n R_j \tag{18}$$

Where R_j are the ranks of sample n and S_n is the sum of the ranks of the sample expected to have the smaller sum

For the Mann-Whitney test, calculate the U statistic for the sample that is expected to have the smaller sum under the alternative hypothesis.

$$U_m$$
 = the sum of the observations in n exceeding each observation in m (19)

$$U_n$$
 = the sum of the observations in m exceeding each observation in n (20)

There is a linear relation between S_n and U_n . It is expressed as

$$U_m = S_m - \frac{1}{2}m(m+1) \tag{21}$$

and similarly,

$$U_n = S_n - \frac{1}{2}n(n+1) \tag{22}$$

where

$$U_m = mn - U_n \tag{23}$$

In a two-tailed test, use the smallest U statistic to test for significance.

Large sample sizes.

The large-sample approximation using the Wilcoxon statistic, S_n is:



$$z = \frac{S_n - \frac{n(n+m+1)}{2}}{\sqrt{\frac{mn(m+n+1)}{12}}}$$
(24)

The large-sample approximation with the U statistic is

$$z = \frac{U + \frac{1}{2} - \frac{1}{2}mn}{\sqrt{\frac{mn(m+n+1)}{12}}}$$
 (25)

In either case, reject H_0 if $z < -z_{\alpha}$ (or $z < -z_{\alpha/2}$ for a two-tailed test).

Example.

The Wilcoxon (Mann-Whitney) Rank Sum statistic was calculated with data from Sample 1 and Sample 5 in Table 3 (Appendix), n = m = 15. The combined samples were ranked, using midranks for ties. The rank sum for Sample 1 was 258.5 and for Sample 5, 206.5. Hence S = 206.5. Calculating the U statistic, U = 206.5-0.5(15)(16) = 86.5. The large sample approximation for the U statistic is $\frac{86.5 + .5 - .5(15^2)}{\sqrt{15^2(31)}} = \frac{-25.5}{24.1091} = -1.05769$. Because -1.05769 > -1.95996, the null hypothesis cannot be

rejected.

Hodges-Lehmann Estimator of the Difference in Location

When a difference in location exists, it may be appropriate to develop an estimate of the difference. Suppose there are two populations that are assumed to have similar shaped distributions, but have different locations. The problem is to develop a confidence interval that will have the probability of $1 - \alpha$ that the actual difference lies in the interval

Procedure.

All the pairwise differences are computed, $x_i - y_j$. For sample sizes of m and n, there are mn differences. The differences are put in ascending order. The task is to find two integers l and u such that the probability that the difference lies between l and u is equal to $1 - \alpha$. These limits are chosen symmetrically. The appropriate lower tail critical value is found for the Mann-Whitney U statistic. This value is the upper limit of the lower end of the differences. Therefore l is the next consecutive integer. The upper limit of the confidence interval is the lth difference from the upper end. Using the relationship



l + u = mn + 1, u = mn - l + 1. The interval (l, u) is the confidence interval for the difference in location for the two populations.

Large sample sizes.

"l and u may be approximated by

$$l = \left[\frac{mn}{2} - z_{\alpha/2} \sqrt{\frac{mn(m+n+1)}{12}} - \frac{1}{2} \right]$$
 (26)

$$u = \left[\frac{mn}{2} + z_{\alpha/2} \sqrt{\frac{mn(m+n+1)}{12}} - \frac{1}{2} \right]$$
 (27)

where the square brackets denote integer nearest to the quantity within, and $z_{\alpha/2}$ is the suitable upper critical point of the standard normal distribution" (Deshpande, Gore, & Shanubhogue, 1995, p. 45).

Example.

The Hodges-Lehmann estimate of the difference in location was computed using Samples 1 and 5 in Table 3 (Appendix), n = m = 15. All possible differences were computed and ranked. Using the large sample approximation formula (26), l = 112.5 - 1.95596 (24.109) - .5 = 64.844. Thus l = 65 and the lower bound is the 65th difference, - 4. The upper bound is the 65th difference from the upper end, or the 225 - 65 + 1 = 161st value, 14. The confidence interval is (-4, 14).

Siegel-Tukey Test

In 1960, the Siegel-Tukey test was developed, which is similar in procedure to the Wilcoxon Rank Sum test for difference in location. This test is based upon the logic that if two samples come from populations with the same median, the one with the greater variability will have more extreme scores. An advantage of the Siegel-Tukey statistic is that it uses the Wilcoxon table of critical values or can be transformed into a U statistic for use with the Mann-Whitney U table of critical values.

The hypotheses for a two-tailed test are:

 H_0 : There is no difference in spread between the two populations

 H_1 : There is some difference in spread between the two populations

Procedure.

The two combined samples are ordered, keeping track of sample membership. The ranking proceeds as follows: the lowest observation is ranked 1, the highest is ranked 2, and the next highest 3. Then the second lowest is ranked 4 and the subsequent observation ranked 5. The ranking continues to



alternate from lowest to highest, ranking two scores at each end. If there is an odd number of scores, the middle score is discarded and the sample size reduced accordingly. Below is an illustration of the ranking procedure.

where N = n + m.

Test statistic

The sum of ranks is calculated for one sample. The rank sum can be used with a table of critical values or it can be transformed into a U statistic by the following formula.

$$U^* = R_n - \frac{1}{2}n(n+1) \tag{28}$$

or

$$U^{\bullet} = R_m - \frac{1}{2}m(m+1) \tag{29}$$

Large sample sizes.

The large-sample approximations are the same for the Siegel-Tukey test as for the Wilcoxon Rank Sum or the Mann-Whitney U statistic, formulas (24) and (25).

Example.

The Siegel-Tukey statistic was calculated using Sample 1 and Sample 5 in Table 3 (Appendix), n = m = 15. The samples were combined and ranked according to the method described. Then tied ranks were averaged. The sum of ranks was 220.5 for Sample 1, and 244.5 for Sample 5. The U statistic is

220.5 - .5(15)(16) = 100.5. The large sample approximation is
$$z = \frac{100.5 + .5 - .5(15^2)}{\sqrt{\frac{15^2(31)}{12}}} = \frac{-11.5}{24.109127} = -\frac{11.5}{24.109127}$$

0.476998. Because -0.476998 > -1.95996, the null hypothesis cannot be rejected.

The Mood Test

In 1954, the Mood test was developed based on the sum of squared deviations of one sample's ranks from the average combined ranks. The null hypothesis is that there is no difference in spread against the alternative hypothesis that there is some difference.

Procedure.

Let sample 1 be $x_1, x_2, ..., x_m$ and let sample 2 be $y_1, y_2, ..., y_n$. Arrange the combined samples in



ascending order and rank the observations from 1 to m + n. Let R_i be the rank of x_i . Let N = m + n. If Nis odd, the middle rank is ignored to preserve symmetry.

Test statistic.

The test statistics is

$$M = \sum_{i=1}^{m} \left(R_i - \frac{m+n+1}{2} \right)^2 . \tag{30}$$

Large sample sizes.

The large sample approximation is

$$z = \frac{M - \frac{m(N^2 - 1)}{12}}{\sqrt{\frac{mn(N+1)(N^2 - 4)}{180}}}$$
(31)

where N = m + n and M is the test statistic.

Example.

The Mood statistic was calculated using Sample 1 and Sample 5 in Table 3 (Appendix), n = m =15. The combined samples are ranked, with midranks assigned to ties. The overall mean of the ranks is 15.5, and the sum of squared deviations of the ranks from the mean for Sample 1 was calculated, yielding M = 1257. The large sample approximation is $\frac{1257 - 1123.75}{\sqrt{34720}} = \frac{133.25}{186.333} = 0.71512$. Because 0.71512 < 1.95596, the null hypothesis cannot be rejected.

The Savage Test for Positive Random Variables

Unlike the Siegel-Tukey test and the Mood test, the Savage test does not assume that location remains the same. It is assumed that differences in scale cause a difference in location. The samples are assumed to be drawn from continuous distributions.

The null hypothesis is that there is no difference in spread against the two-tailed alternative, there is a difference.

Procedure.

Let sample 1 be $x_1, x_2, ..., x_m$ and let sample 2 be $y_1, y_2, ..., y_n$. The combined samples are ordered, keeping track of sample membership. Let R_i be the rank for x_i . The test statistic is computed for either sample.



Test statistic.

The test statistic is

$$S = \sum_{i=1}^{m} a(R_i) \tag{32}$$

where

$$a(i) = \sum_{j=N+1-i}^{N} \frac{1}{j}$$
 (33)

such that
$$a(1) = \frac{1}{N}$$
, $a(2) = \frac{1}{N-1} + \frac{1}{N}$, ..., $a(N) = 1 + \frac{1}{2} + \frac{1}{3} + ... + \frac{1}{N-1} + \frac{1}{N}$.

Large sample sizes.

For large sample sizes the following normal approximation may be used.

$$S^* = \frac{S - n}{\sqrt{\frac{nm}{N - 1} \left(1 - \frac{1}{N} \sum_{j=1}^{N} \frac{1}{j}\right)}}$$
(35)

 S^{\bullet} is compared to the critical z value from the standard normal distribution.

Example.

The Savage statistic was calculated using Sample 1 and Sample 5 in Table 3 (Appendix), n = m = 15. Using Sample 1, S = 18.3114. The large sample approximation is $\frac{18.3114 - 15}{\sqrt{7.7586(.86683)}} = \frac{3.114}{2.59334} = 1.27689$. Because 1.27689 < 1.95596, the null hypothesis cannot be rejected.

Ansari-Bradley Test

This is a rank test for spread when the population medians are the same. The null hypothesis is that the two populations have the same spread against the two-tailed alternative that the spreads of the two populations differ.

Procedure.

Order the combined samples, keeping track of sample membership. Rank the smallest and largest observation 1. Rank the second lowest and second highest 2. If the combined sample size, N, is odd, the middle score will be ranked $\frac{N+1}{2}$ and if N is even the middle two ranks will be $\frac{N}{2}$. The pattern will be



either 1, 2, 3, ...,
$$\frac{N+1}{2}$$
, ..., 3, 2, 1 (N odd), or 1, 2, 3, ..., $\frac{N}{2}$, $\frac{N}{2}$, ..., 3, 2, 1 (N even).

Test statistic.

The test statistic, W, is the sum of the ranks of sample 1.

$$W = \sum_{i=1}^{m} R_i \tag{35}$$

where R_i is the rank of the i^{th} observation of a sample.

Large sample sizes.

There are two formulae, one if N is even, and one if N is odd.

$$W^* = \frac{W - \frac{m(m+n+2)}{4}}{\sqrt{\frac{mn(m+n+2)(m+n-2)}{48(m+n-1)}}}$$
(36)

if N is even and

$$W^* = \frac{W - \frac{m(m+n+1)^2}{4(m+n)}}{\sqrt{\frac{mn(m+n+1)[3+(m+n)^2]}{48(m+n)^2}}}$$
(37)

if N is odd. Reject the null hypothesis if $W \ge z_{\alpha/2}$.

Example.

The Ansari-Bradley statistic was calculated using Sample 1 and Sample 5 in Table 3 (Appendix), n=m=15. The combined samples were ranked using the method described, and tied ranks were assigned average ranks. The statistic, W, is 126.5, the rank sum of Sample 5. The large sample approximation is $\frac{126.5-120}{\sqrt{144.827586}} = \frac{6.5}{12.03443} = 0.540117$. Because 0.540117 < 1.95596, the null hypothesis cannot be rejected.

Comparisons of Several Populations

This section considered tests against an omnibus alternative and tests involving an ordered hypothesis. The omnibus tests were the Kruskal-Wallis test and Friedman's test. The tests for ordered alternatives are the Terpstra-Jonckheere test, Page's test, and the match test.

The Kruskal-Wallis test is a test for independent samples. It is analogous to the one-way analysis

of variance. Friedman's test is an omnibus test for k related samples, and is analogous to a two-way analysis of variance.

Comparisons of several populations with ordered alternative hypotheses are extensions of a one-sided test. Where an omnibus alternative states only that there is some difference between the populations, an ordered alternative specifies the order of differences. Three tests for an ordered alternative were included, the Terpstra-Jonckheere Test, Page's Test, and the Match Test.

Kruskal-Wallis Test

In 1952, the Kruskal-Wallis test was derived from the F test. It is an extension of the Wilcoxon (Mann-Whitney) test. The null hypothesis is that the k populations have the same average (median). The alternative hypothesis is that at least one sample is from a distribution with a different average (median).

Procedure.

Rank all the observations in the combined samples, keeping track of the sample membership. Compute the rank sums of each sample. Let R_i equal the sum of the ranks of the ith sample of sample size n_i . The logic of the test is that the ranks should be randomly distributed among the k samples.

Test statistic.

The formula is

$$H = \frac{12}{N(N+1)} \sum_{i=1}^{k} \frac{R_i^2}{n_i} - 3(N+1)$$
 (38)

where N is the total sample size, n_i is the size of the ith group, k is the number of groups, and R_i is the rank-sum of the ith group. Reject H_0 when $H \ge \text{critical value}$.

Large sample sizes.

For large sample sizes, the null distribution is approximated by the χ^2 distribution with k-1 degrees of freedom. Thus, the rejection rule is to reject H_0 if $H \ge \chi^2_{a,k-1}$ where $\chi^2_{a,k-1}$ is the value of χ^2 at nominal α with k-1 degrees of freedom.

Example.

The Kruskal-Wallis statistic was calculated using Samples 1 – 5 in Table 3 (Appendix), $n_1 = n_2 = n_3 = n_4 = n_5 = 15$. The combined samples were ranked, and tied ranks were assigned midranks. The rank sums were: $R_1 = 638$, $R_2 = 595$, $R_3 = 441.5$, $R_4 = 656.5$, and $R_5 = 519$. The sum of $R_i^2 = 1,656,344.5$, i = 1,656,344.5



1, 2, 3, 4, 5.
$$H = \frac{12}{75.76} \left(\frac{1,656,344.5}{15} \right) - 3.76 = 0.00211(110,422.9667) - 228 = 4.4694$$
. The statistic, $H = \frac{12}{15} \left(\frac{1,656,344.5}{15} \right) - \frac{1}{15} \left$

4.4694. The large sample approximation is chi-square with 5-1=4 degrees of freedom at $\alpha=.05$ which is 9.488. Because 4.4694 < 9.488, the null hypothesis cannot be rejected.

Friedman's Test

In 1937, the Friedman test was developed as a test for k related samples. The null hypothesis is that the samples come from the same population against the alternative that at least one of the samples comes from a different population. The data are arranged in k columns and n rows, where each row contains k related observations.

Procedure.

Rank the observations for each row from 1 to k. For each of the k columns, the ranks are added and averaged, and the mean is designated \overline{R}_j . The overall mean of the ranks is $\overline{R} = \frac{1}{2}(k+1)$. The sum of the squares of the deviations of mean of the ranks of the columns from the overall mean rank is computed. The test statistic is a multiple of this sum.

Test statistic.

The test statistic for Friedman's test is M, which is a multiple of S, as follows:

$$S = \sum_{j=1}^{k} (\overline{R}_j - \overline{R})^2 \tag{39}$$

$$M = \frac{12n}{k(k+1)}S\tag{40}$$

where n is the number of rows, and k is the number of columns. An alternate formula that does not use S is as follows.

$$M = \frac{12}{nk(k+1)} \sum_{j=1}^{k} R_j^2 - 3n(k+1)$$
 (41)

where *n* is the number of rows, *k* is the number of columns, and R_j is the rank sum for the j^{th} column, j = 1, 2, 3, ..., k.

Large sample sizes.

For large sample sizes, the critical values can be approximated by the chi-square distribution with k-1 degrees of freedom.

Example.

Friedman's statistic was calculated with Samples 1-5 in Table 3 (Appendix), $n_1=n_2=n_3=n_4=n_5=15$. The rows were ranked, with midranks assigned to tied ranks. The column sums are: $R_1=48.5$, $R_2=47$, $R_3=33$, $R_4=52.5$, and $R_5=44$. The sum of the squared rank sums is 10,342.5. The statistic, $M=\frac{12}{15\cdot 5\cdot 6}(10,342.5)-3\cdot 15\cdot 6=0.02667(10,342.5)-270=5.8$. The large sample approximation is chi-square with 5-1=4 degrees of freedom and $\alpha=.05$ which is 9.488. Because 5.8<9.488, the null hypothesis cannot be rejected.

Terpstra-Jonckheere Test

This is a test for more than two independent samples. It was first developed by Terpstra in 1952 and later independently developed by Jonckheere in 1954. The null hypothesis is that the medians of the samples are equal against the alternative that the medians are either decreasing or increasing. This test is based on the Mann-Whitney U statistic, where U is calculated for each pair of samples and the U statistics are added.

Suppose the null hypothesis is $H_0: m_1 = m_2 = ... = m_k$ and the alternative hypothesis is $H_1: m_1 < m_2 < ... < m_k$ for i = 1, 2, ..., k, where m_i is the median for sample i. The U statistic is calculated for each of the $\frac{k(k-1)}{2}$ pairs, which are ordered so that the smallest U is calculated.

Test statistic.

The test statistic is the sum of the U statistics.

$$W = U_{k1} + U_{k2} + \dots + U_{31} + U_{32} + U_{21}$$
 (42)

where U_{ij} is the number of (x_i, x_j) pairs with x_j less than x_i

Large sample sizes.

The null distribution of W approaches normality as the sample size increases. The mean of the distribution is

$$\mu = \frac{(N^2 - \sum n_i^2)}{4} \tag{43}$$

and the standard deviation is

$$\sigma = \sqrt{\frac{N^2(2N+3) - \sum n_i^2(2n_i + 3)}{72}}$$
 (44)



BEST COPY AVAILABLE

The critical value for large samples is given by

$$W \le \mu - z\sigma - \frac{1}{2} \tag{45}$$

where z is the standard normal value, and $\frac{1}{2}$ is a continuity correction.

Example.

The Terpstra-Jonckheere statistic was calculated with Samples 1 – 5 in Table 3 (Appendix), $n_1 = n_2 = n_3 = n_4 = n_5 = 15$. This was done as a one-tailed test with $\alpha = .05$. The U statistics for each sample were calculated. $U_{2,1} = 121$, $U_{3,1} = 145$, $U_{4,1} = 103$, $U_{5,1} = 135$, $U_{3,2} = 142$, $U_{4,2} = 97$, $U_{5,2} = 124$, $U_{4,3} = 71$, $U_{5,3} = 91$, and $U_{5,4} = 136$, for a total W = 1165. The large sample approximation was calculated, with $\mu = 1125$ and $\sigma = 106.94625$. The approximation is 1125 - 1.6449(106.9463) - .5 = 948.584. Because 1165 > 948.584 the null hypothesis cannot be rejected.

Page's Test

In 1963, Page's test for an ordered hypothesis for k > 2 related samples was developed. It takes the form of a randomized block design, with k columns and n rows. The null hypothesis is $H_0: m_1 = m_2 = \ldots = m_k$ and the alternative hypothesis is $H_1: m_1 < m_2 < \ldots < m_k$ for $i = 1, 2, \ldots k$. For this test, the alternative must be of this form. The samples need to be reordered if necessary.

Procedure.

The data are ranked from 1 to k for each row, creating a table of the ranks. The ranks of each of the k columns are totaled. If the null hypothesis is true, the ranks should be evenly distributed over the columns, whereas if the alternative is true, the ranks sums should increase with the column index.

Test statistic.

Each column rank-sum is multiplied by the column index. The test statistic is

$$L = \sum_{i=1}^{k} iR_i \tag{46}$$

where i is the column index, i = 1, 2, 3, ..., k, and R_i is the rank sum for the i^{th} column.

Large sample sizes

The mean of L is

$$\mu = \frac{nk(k+1)^2}{4} \tag{47}$$



and the standard deviation is

$$\sigma = \sqrt{\frac{nk^2(k+1)(k^2-1)}{144}} \tag{48}$$

For a given α , the approximate critical region is

$$L \ge \mu + z\sigma + \frac{1}{2} \tag{49}$$

Example.

Page's statistic was calculated with Samples 1 – 5 in Table 3 (Appendix), $n_1 = n_2 = n_3 = n_4 = n_5 = 15$. This was done as a one-tailed test with $\alpha = .05$. The rows are ranked with midranks assigned to tied ranks. The column sums are: $R_1 = 48.5$, $R_2 = 47$, $R_3 = 33$, $R_4 = 52.5$, and $R_5 = 44$. The statistic, L, is the sum of $iR_i^2 = 671.5$, where i = 1, 2, 3, 4, 5. The large sample approximation was calculated with $\mu = 675$ and $\sigma = 19.3649$. The approximation is 675 + 1.64485(19.3649) + .5 = 707.352. Because 671.5 < 707.352, the null hypothesis cannot be rejected.

The Match Test for Ordered Alternatives

The match test is a test for k > 2 related samples with an ordered alternative hypothesis. The match test was developed by Neave and Worthington (1988). It is very similar in concept to Page's test, but instead of using rank-sums, it uses the number of matches of the ranks with the expected ranks plus half the near matches.

The hypotheses are the same as for Page's test. The null hypothesis is $H_0: m_1 = m_2 = \dots = m_k$ and the alternative hypothesis is $H_1: m_1 < m_2 < \dots < m_k$ for $i = 1, 2, \dots k$.

Procedure.

A table of ranks is compiled with the observations in each row ranked from 1 to k. Ties are assigned average ranks. Each rank, r_i , is compared with the expected rank, i, which is the column index. If the rank equals the column index, it is a match. The number of matches is counted. Every non-match such that $0.5 \le |r_i - i| \le 1.5$ is counted as a near match.

Test statistic.

The test statistic is

$$L_2 = L_1 + \frac{1}{2} \text{(number of near matches)}$$
 (50)

where L_1 is the number of matches.



Large sample sizes.

The null distribution approaches a normal distribution for large sample size. The mean and standard deviation for L_2 are as follows:

$$\mu = n\left(2 - \frac{1}{k}\right) \tag{51}$$

$$\sigma = \sqrt{\frac{n}{k} \left(\frac{3(k-2)}{2}\right) + \frac{1}{k(k-1)}} \tag{52}$$

For a given level of significance α the critical value approximation is

$$L_2 \ge \mu + z\sigma + \frac{1}{2} \tag{53}$$

where z is the upper-tail critical value from the standard normal distribution and $\frac{1}{2}$ is a continuity correction.

Example.

The Match statistic was calculated with Samples 1 – 5 in Table 3 (Appendix), $n_1 = n_2 = n_3 = n_4 = n_5 = 15$. This was done as a one-tailed test with $\alpha = .05$. The rows are ranked with midranks assigned to tied ranks. The number of matches for the five columns are 3, 3, 2, 2, and 1, for $L_1 = 11$. The number of near matches were 1, 6, 8, 8, and 4, for $L_2 = 27$. The statistic, L = 11 + .5(27) = 24.5. For the large sample approximation, $\mu = 27$ and $\sigma = 3.68103$. The approximation is 27 + 1.6449(3.68103) + .5 = 33.5549. Because 24.5 < 33.5549, the null hypothesis cannot be rejected.

Rank Correlation Tests

The rank correlation is a measure of the association of a pair of variables. Two tests of association were studied, Spearman's rank correlation coefficient (rho) and Kendall's rank correlation coefficient (tau).

Spearman's Rank Correlation Coefficient

Spearman's rank correlation (rho) was published in 1904. Let X and Y be the two variables of interests. Each observed pair is denoted by (x_i, y_i) . The paired ranks are denoted by (r_i, s_i) where r_i is the rank of x_i and s_i is the rank of y_i . The null hypothesis for a two-tailed test is $H_0: \rho = 0$ against the



alternative, $H_1: \rho \neq 0$. The alternative hypotheses for a one-tailed test are $H_1: \rho > 0$ or $H_1: \rho < 0$.

Procedure.

Rank both X and Y scores while keeping track of the original pairs. Form the rank pairs (r_i, s_i) which correspond to the original pair, (x_i, y_i) . Calculate the sum of the squared differences between r_i and s_i .

Test statistic.

If there are no ties, the formula is

$$\rho = 1 - \frac{6T}{n(n^2 - 1)} \tag{54}$$

where

$$T = \sum (r_i - s_i)^2 \tag{55}$$

Large sample sizes.

For large n the distribution of ρ is approximately normal. The critical values can be found by $z = \rho \sqrt{n-1}$. The rejection rule for a two-tailed test is to reject H_0 if $z > z_{\alpha/2}$ or $z < -z_{\alpha/2}$ where $z_{\alpha/2}$ is the critical value for the given level of significance.

Example.

Spearman's rho was calculated using Sample 1 and Sample 5 in Table 3 (Appendix), n = 15. The sum of the squared rank differences for the two samples is T = 839. Rho is $1 - \frac{6(839)}{15(224)} = 1 - \frac{5034}{3360} = 1 - 1.498214 = -0.498214$. The large sample approximation is z = -0.498214 $\sqrt{14} = -1.864147$. Because -1.864 > -1.956, the null hypothesis cannot be rejected.

Kendall's Rank Correlation Coefficient

Kendall's rank correlation coefficient (tau) is similar to Spearman's rho. The underlying concept is the tendency for concordance. Concordance is the concept that if $x_i > x_j$ then $y_i > y_j$. Concordance implies that the differences $x_i - x_j$ and $y_i - y_j$ have the same sign, either "+" or "-". Discordant pairs are pairs that have opposite signs, that is, $x_i > x_j$ but $y_i < y_j$, or the opposite, $x_i < x_j$ but $y_i > y_j$. A high number of concordant pairs support the alternative hypothesis of positive, and correlation, a high number of discordant pairs support an alternative hypothesis of negative correlation.



Procedure.

Arrange the pairs in ascending order of X. Count the number of y_i which are smaller than y_i . This is the number of discordant pairs (N_D) for x_i . Repeat the process for each subsequent x_i , counting the number of smaller y_j to the right of the y_i , j = i + 1, i + 2, i + 3, ..., n.

Test statistic.

Because the total number of pairs is $\frac{1}{2}n(n-1)$, it follows that $N_c = \frac{1}{2}n(n-1) - N_D$. The statistic (τ) is defined as

$$\tau = \frac{N_C - N_D}{\frac{1}{2}n(n-1)} \tag{56}$$

This formula can be simplified by substituting $N_c = \frac{1}{2}n(n-1) - N_D$ into the formula so that

$$\tau = 1 - \frac{4N_D}{n(n-1)} \tag{57}$$

From this formula, it can be seen that if there are no discordant pairs, τ equals 1, showing positive correlation. If all pairs are discordant, $4N_D = 4(\frac{1}{2})n(n-1) = 2n(n-1)$ and $\tau = 1 - 2 = -1$, showing negative correlation.

Large sample sizes.

For large sample sizes, the formula is

$$z = \frac{3\tau\sqrt{n(n-1)}}{\sqrt{2(2n+5)}}$$
 (58)

where z is compared to the z score from the standard normal distribution for the appropriate alpha level.

Example.

Kendall's tau was calculated using Sample 1 and Sample 5 in Table 3 (Appendix), n = 15. The number of discordant pairs for each pair, (x_1, x_5) , were 12, 8, 8, 5, 9, 5, 6, 3, 5, 3, 0, 3, 0, 1, and 0. The total number of discordant pairs, N_D is 68. Tau is $1 - \frac{4 \cdot 68}{15 \cdot 14} = 1 - \frac{272}{210} = -0.295238$. The large sample approximation is $\frac{3(-.295238)\sqrt{(15)(14)}}{\sqrt{2(35)}} = \frac{-12.83522}{8.3666} = -1.534102$. Because -1.534102 > -1.95596, the

null hypothesis cannot be rejected.

References

- *Anderson, D. R., Sweeney, D. J., & Williams, T. A. (1999). <u>Statistics for Business and Economics</u> (7th ed.). Cincinnati: South-Western College Publishing Co.
- *Berenson, M. L., Levine, D. M., & Rindskopf, D. (1988). <u>Applied statistics: A first course.</u> Englewood Cliffs, NJ: Prentice Hall, Inc.
- Blair, R. C., & Higgins, J. J. (1985). A comparison of the power of the paired samples rank transformation to that of Wilcoxon's signed rank statistic. <u>Journal of Educational Statistics</u>, 10, 368-383.
- Blair, R. C., & Higgins, J. J. (1980). A comparison of the power of Wilcoxon's rank-sum statistic to that of Student's t statistic under various nonnormal distributions. <u>Journal of Educational Statistics</u>, 5, 309-335.
- Blair, R. C., Sawilowsky, S. S., & Higgins, J. J. (1987). Limitations of the rank transformation in factorial ANOVA. <u>Communications in Statistics</u>, 16, 1133-1145.
 - Bradley, J. V. (1968). Distribution-free statistical tests. Englewood Cliffs, NJ: Prentice-Hall Inc.
- Bradley, J. V. (1978). Robustness? <u>British Journal of Mathematical and Statistical Psychology</u>, <u>31</u>, 144 152.
- Bradstreet, T. E. (1997). A Monte Carlo study of Type I error rates for the two-sample Behrens-Fisher problem with and without rank transformation. <u>Computational Statistics and Data Analysis</u>, <u>25</u>, 167-179.
- Bridge, P. K., & Sawilowsky, S. S. (1999). Increasing physician's awareness of the impact of statistical tests on research outcomes: Investigating the comparative power of the Wilcoxon Rank-Sum test and independent samples t test to violations from normality. <u>Journal of Clinical Epidemiology</u>, <u>52</u>, 229-235.
 - Conover, W. J. (1971). Practical nonparametric statistics. New York: John Wiley & Sons, Inc.
- *Daly, F., Hand, D. J., Jones, M. C., Lunn, A. D., & McConway, K. J. (1995). <u>Elements of statistics.</u> Workingham, England: Addison-Wesley.
 - *Daniel, W. W. (1978). Applied nonparametric statistics. Boston: Houghton Mifflin Co.
- Deshpande, J.V., Gore, A. P., & Shanubhogue, A. (1995). <u>Statistical analysis of nonnormal data.</u> New York: John Wiley & Sons, Inc.
- *Ferguson, G. A. (1971). <u>Statistical analysis in psychology and education</u> (3rd ed.). New York: McGraw Hill book Company.
 - *Ferguson, G. A. (1981). Statistical analysis in psychology and education (5th ed.). New York:



BEST COPY AVAILABLE

McGraw - Hill Book Company.

*Gravetter, F. J., & Wallnau, L. B. (1985). <u>Statistics for the behavioral sciences.</u> St. Paul: West Publishing Co.

Gibbons, J. D. (1971). Nonparametric statistical inference. New York: McGraw-Hill Book Company.

Hájek, J. (1969). A course in nonparametric statistics. San Francisco: Holden - Day.

Harwell, M., & Serlin, R. C. (1997). An empirical study of five multivariate tests for the single-factor repeated measures model. <u>Communications in Statistics</u>, 26, 605-618.

*Hays, W. L. (1994). Statistics (5th ed.). Fort Worth: Harcourt Brace College Publishers.

Headrick, T. C., & Sawilowsky, S. S. (2000). Type I error and power of the RT ANCOVA. American Educational Research Association, SIG/Educational Statisticians. New Orleans, LA

Headrick, T. C., & Sawilowsky, S. S. (1999). Type I error and power of the rank transform in factorial ANCOVA. Statistics Symposium on Selected Topics in Nonparametric Statistics. Gainesville, FL.

*Hildebrand, D. (1986). <u>Statistical thinking for behavioral scientists.</u> Boston: Duxbury Press. Hollander, M. & Wolfe, D. (1973). <u>Nonparametric statistical methods.</u> New York: John Wiley & Sons.

*Jarrett, J. & Kraft, A. (1989). <u>Statistical analysis for decision making</u>. Boston: Allyn and Bacon. Kelley, D. L., & Sawilowsky, S. S. (1997). Nonparametric alternatives to the F statistic in analysis of variance. <u>Journal of Statistical Computation and Simulation</u>, 58, 343-359.

*Knoke, D. and Bohrnstedt, G. W. (1991). <u>Basic social statistics</u>. New York: F. E. Peacock Publishers, Inc.

*Kraft, C. H. & van Eeden, C. (1968). A nonparametric introduction to statistics. New York: Macmillan Co.

*Krauth, J. (1988). <u>Distribution-free statistics: An application-oriented approach</u>. Amsterdam: Elsevier.

*Kurtz, N. R. (1983). <u>Introduction to social statistics</u>. New York: McGraw - Hill Book Company.

Lahey (1998). Essential Lahey Fortran 90. Incline Village, NY: Lahey Computer Systems, Inc.

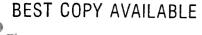
Laubscher, N. F., Steffens, F. E., & De Lange, E. M. (1968). Exact critical values for Mood's distribution-free test statistic for dispersion and its normal approximation. <u>Technometrics</u>, 10, 497-507.



*Lehmann, E. L. & D'Abrera, H. J. M. (1975). Nonparametric statistical methods based on ranks. New York: McGraw-Hill International Book Company.

Ludbrook, J. L., & Dudley, H. (1998). Why permutation tests are superior to t and F tests in biomedical research. The American Statistician, 52, 127-132.

- *Manoukian, E. B. (1986). <u>Mathematical nonparametric statistics.</u> New York: Gordon & Breach Science Publications.
- *McClave, J. T., Dietrich II, F. H. (1988). <u>Statistics</u> (4th ed.). San Francisco: Dellen Publishing Company.
- *Mendenhall, W. & Reinmuth, J. E. (1978). <u>Statistics for management and economics</u> (3rd ed.). North Scituate, MA: Duxbury Press.
- Micceri, T. (1989). The unicorn, the normal curve, and other improbable creatures. Psychological Bulletin, 105, 156-166.
- *Montgomery, D. C., & Runger, G. C. (1994). <u>Applied statistics and probability for engineers</u>. New York: John Wiley and Sons, Inc.
- Musial, J., III. (1999). Comparing exact tests and asymptotic tests with colorectal cancer variables within the National Health and Nutrition Examination Survey III. Unpublished doctoral dissertation, Wayne State University, Detroit, MI.
- Nanna, M. J. (1997). Robustness and comparative power properties of Hotelling's T² versus the rank transformation test using real pre-test/post-test likert scale data. Unpublished doctoral dissertation, Wayne State University, Detroit, MI.
- Nanna, M. J., & Sawilowsky, S. S. (1998). Analysis of likert scale data in disability and medical rehabilitation evaluation. <u>Psychological Methods</u>, <u>3</u>, 55-67.
 - Neave, H. R., & Worthington, P. L. (1988). Distribution-free tests. London: Unwin Hyman Ltd.
- *Newmark, J. (1988). <u>Statistics and probability in modern life</u> (4th ed.). New York: Saunders College Publishing.
- Posch, M.A., & Sawilowsky, S. (1997). A comparison of exact tests for the analysis of sparse contingency tables. Joint Statistical Meetings, American Statistical Association. Anaheim, CA.
- *Rosenberg, K. M. (1990). <u>Statistics for behavioral scientists</u>. Dubuque, IA: Wm. C. Brown Publishers.
- *Runyon, R. P. (1977). <u>Nonparametric statistics: A contemporary approach.</u> Reading MA: Addison-Wesley Publishing Co.



Sawilowsky, S. S. (1985). Robust and power analysis for the 2x2x2 ANOVA, rank transformation, random normal scores, and expected normal scores transformation tests. Unpublished doctoral dissertation, University of South Florida, Tampa, FL.

Sawilowsky, S. S. (1989). Rank transformation: the bridge is falling down. American Educational Research Association, SIG/Educational Statisticians, San Francisco, CA.

Sawilowsky, S. S. & Blair, R. C. (1992). A more realistic look at the robustness and Type II error properties of the *t* test to departures from population normality. <u>Psychological Bulletin, 111,</u> 352-360.

Sawilowsky, S. S., & Brown, M. T. (1991). On using the t test on ranks as an alternative to the Wilcoxon test. <u>Perceptual and Motor Skills</u>, 72, 860-862.

Sawilowsky, S. S., Blair, R. C., & Higgins, J. J. (1989). An investigation of the type I error and power properties of the rank transform procedure in factorial ANOVA. <u>Journal of Educational Statistics</u>, <u>14</u>, 255-267.

Sawilowsky, S. S., Blair, R. C., & Micceri, T. (1990). REALPOPS.LIB: A PC FORTRAN library of eight real distributions in psychology and education. <u>Psychometrika</u>, <u>55</u>, 729.

Siegel, S. & Castellan, Jr., N. J. (1988). <u>Nonparametric statistics for the behavioral sciences.</u>
New York: McGraw-Hill, Inc.

*Snedecor, G. W. & Cochran, W. G. (1967). Statistical methods. Ames, IA: Iowa State University Press.

Sprent, P. (1989). Applied nonparametric statistical methods. London: Chapman and Hall.

*Triola, M. (1995). <u>Elementary statistics</u> (6th ed.). Reading MA: Addison – Wesley Publishing Company.

*Wilcox, R. R. (1996). <u>Statistics for the social sciences.</u> San Diego: Academic Press.

*Zikmund, W. G. (1991). Business research methods (3rd ed.). Chicago: The Dryden Press.



^{*} References marked by the "*" symbol were surveyed for the literature review.

Table 3. Samples Randomly Selected from Micerri's Multimodal Lumpy Data Set.

Sample 1	Sample 2	Sample 3	Sample 4	Sample 5
20	11	9	34	10
33	34	14	10	2
4 .	23	33	38	32
34	37	5	41	4
13	11	8	4	33
6	24	14	26	19
29	5	20	10	11
17	9	18	21	21
39 .	11	8	13	9
26	33	22	15	31
13	32	11 -	35	12
9	18	33	43	20
33	27	20	13	33
16	21	7	20	15
36	8	7	13	15



Table 4. Multimodal Lumpy Set (Micceri, 1989).

Score	Cumulative Frequency	cdf	Score	Cumulative Frequency	cdf
. 0	5	0.01071	22	269	0.57602
1	13	0.02784	23	279	0.59743
2	21	0.04497	24	282	0.60385
3	24	0.05139	25	287	0.61456
4	32 ·	0.06852	26	297	0.63597
5	38	0.08137	27	306	0.65525
6	41	0.08779	28	309	0.66167
7	50	0.10707	29	319	0.68308
8	62	0.13276	30	325	0.69593
9	80	0.17131	31	336	0.71949
10	91	0.19486	32	351	0.75161
11	114	0.24411	33	364	0.77944
12	136	0,29122	34	379	0.81156
13	160	0.34261	35	389	0.83298
14	180	0.38544	36	401	0.85867
15	195	0.41756	37	418	0.89507
16	213	0.45610	38	428	0.91649
17	225	0.48180	39	434	0.92934
18	234	0.50107	40	445	0.95289
19	244	0.52248	41	454	0.97216
20	254	0.54390	42	460	0.98501
21	261	0.55889	43	467	1.00000





(over)



U.S. Department of Education

Office of Educational Research and Improvement (OERI)
National Library of Education (NLE)
Educational Resources Information Center (ERIC)



TM030830

REPRODUCTION RELEASE

(Specific Document)

I. DOCUMENT IDENTIFICATION	:	1 6 1
Title: Review of Tw	enty Nongarametr	ic Statistics
Author(s): Gail Fahoome	: Shlowd S. Saw	i lowsky
Corporate Source:		Publication Date:
II. REPRODUCTION RELEASE:	•	
monthly abstract journal of the ERIC system, Res	timely and significant materials of interest to the eduction of the interest to the eduction (RIE), are usually made available Document Reproduction Service (EDRS). Credit is no notices is affixed to the document.	e to users in microfiche, reproduced paper copy,
If permission is granted to reproduce and disser of the page.	ninate the identified document, please CHECK ONE or	f the following three options and sign at the bottom
The sample sticker shown below will be affixed to all Level 1 documents	The sample sticker shown below will be affixed to all Level 2A documents	The sample sticker shown below will be affixed to all Level 2B documents
PERMISSION TO REPRODUCE AND DISSEMINATE THIS MATERIAL HAS BEEN GRANTED BY	PERMISSION TO REPRODUCE AND DISSEMINATE THIS MATERIAL IN MICROFICHE, AND IN ELECTRONIC MEDIA FOR ERIC COLLECTION SUBSCRIBERS ONLY, HAS BEEN GRANTED BY	PERMISSION TO REPRODUCE AND DISSEMINATE THIS MATERIAL IN MICROFICHE ONLY HAS BEEN GRANTED BY
sample	Sample	Sample
TO THE EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)	TO THE EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)	TO THE EDUCATIONAL RESOURCES INFORMATION CENTER (ERIC)
1	2A	2B
Level 1	Level 2A ↑	Level 2B ↑
Check here for Level 1 release, permitting reproduction and dissemination in microfiche or other ERIC archival media (e.g., electronic) and paper copy.	Check here for Level 2A release, permitting reproduction and dissemination in microfiche and in electronic media for ERIC archival collection subscribers only	Check here for Level 2B release, permitting reproduction and dissemination in microfiche only
	nts will be processed as indicated provided reproduction quality per roduce is granted, but no box is checked, documents will be proce	
as indicated above. Reproduction from the contractors requires permission from the to satisfy information needs of educated.	<u> </u>	ns other than ERIC employees and its system production by libraries and other service agencies
Sign Signatures The Signatures of the Signature of the Signatures of the Signature of the Si	Printed Name/Pos	sition/Title:

48205

III. DOCUMENT AVAILABILITY INFORMATION (FROM NON-ERIC SOURCE):

If permission to reproduce is not granted to ERIC, or, if you wish ERIC to cite the availability of the document from another source, please provide the following information regarding the availability of the document. (ERIC will not announce a document unless it is publicly available, and a dependable source can be specified. Contributors should also be aware that ERIC selection criteria are significantly more stringent for documents that cannot be made available through EDRS.)

	•		
		·	
•		<u> </u>	<u> </u>
			·
	production release is he	OF ERIC TO COPYRIGHT/REPROPRICE TO COPYRIGHT/REPROPRICE TO COPYRIGHT/REPROPRICE TO THE PROPRICE TO THE PROPRIC	OF ERIC TO COPYRIGHT/REPRODUCTION RIGH production release is held by someone other than the addressee, please pro

V. WHERE TO SEND THIS FORM:

Send this form to the following ERIC Clearinghouse:

ERIC CLEARINGHOUSE ON ASSESSMENT AND EVALUATION
UNIVERSITY OF MARYLAND
1129 SHRIVER LAB
COLLEGE PARK, MD 20772
ATTN: ACQUISITIONS

However, if solicited by the ERIC Facility, or if making an unsolicited contribution to ERIC, return this form (and the document being contributed) to:

ERIC Processing and Reference Facility

4483-A Forbes Boulevard Lanham, Maryland 20706

Telephone: 301-552-4200 Toll Free: 800-799-3742 FAX: 301-552-4700 e-mail: ericfac@inet.ed.gov

e-mail: ericfac@inet.ed.gov WWW: http://ericfac.piccard.csc.com

